



(12) **EUROPEAN PATENT APPLICATION**
published in accordance with Art. 153(4) EPC

(43) Date of publication:
18.12.2024 Bulletin 2024/51

(51) International Patent Classification (IPC):
G06N 5/02 (2023.01)

(21) Application number: **22925944.5**

(52) Cooperative Patent Classification (CPC):
G06N 5/02

(22) Date of filing: **10.02.2022**

(86) International application number:
PCT/JP2022/005497

(87) International publication number:
WO 2023/152923 (17.08.2023 Gazette 2023/33)

(84) Designated Contracting States:
AL AT BE BG CH CY CZ DE DK EE ES FI FR GB GR HR HU IE IS IT LI LT LU LV MC MK MT NL NO PL PT RO RS SE SI SK SM TR
Designated Extension States:
BA ME
Designated Validation States:
KH MA MD TN

(72) Inventors:
• **IWASHITA, Hiroaki**
Kawasaki-shi, Kanagawa 211-8588 (JP)
• **ASAI, Tatsuya**
Kawasaki-shi, Kanagawa 211-8588 (JP)
• **UEMURA, Kento**
Kawasaki-shi, Kanagawa 211-8588 (JP)
• **KOYANAGI, Yusuke**
Kawasaki-shi, Kanagawa 211-8588 (JP)

(71) Applicant: **FUJITSU LIMITED**
Kawasaki-shi, Kanagawa 211-8588 (JP)

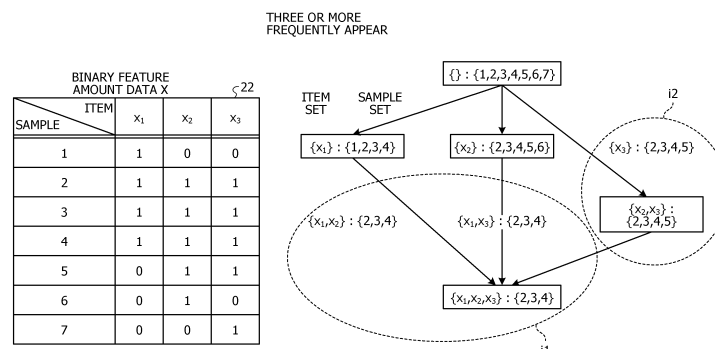
(74) Representative: **Hoffmann Eitle**
Patent- und Rechtsanwälte PartmbB
Arabellastraße 30
81925 München (DE)

(54) **INFORMATION PROCESSING PROGRAM, INFORMATION PROCESSING DEVICE, AND INFORMATION PROCESSING METHOD**

(57) An information processing device (1) generates, from feature amount data (21) in which values of a plurality of feature amounts included in each sample are accumulated for each sample, binary feature amount data (22) obtained by binarizing, for each sample, the values of the plurality of feature amounts included in each sample based on an item set in advance, enumerates, by using the binary feature amount data (22), item sets in

which all sample sets indicate true values, computes, for each item set, a correlation between the plurality of feature amounts in the feature amount data (21) in a sample set associated with each item set, and selects an item set determined to have a correlation as a condition to be causally searched. As a result, it is possible to accurately select a condition under which a correlation appears.

FIG. 3



Description

CITATION LIST

TECHNICAL FIELD

NON-PATENT DOCUMENT

[0001] The present invention relates to an information processing program and the like.

BACKGROUND ART

[0002] In recent years, research has been conducted to efficiently narrow down the number of conditions to be causally searched, by extracting correlated conditions. FIG. 8 is a reference diagram illustrating condition extraction for statistical causal search. As illustrated in FIG. 8, such a technique extracts all condition candidates for correlated data from past data for AI to train. Then, in the technique, all conditions for data having a causal relationship are extracted from the extracted condition candidates for the data. However, in the technique, while all condition candidates for the data are searched for a causal relationship, there is a problem that the search is unrealistic from a viewpoint of an amount of computation.

[0003] Thus, a technique that efficiently narrows down the number of conditions to be causally searched, by relaxing a condition search target to a correlation from the causal relationship has been disclosed (for example, Non-Patent Document). FIG. 9 is a reference diagram illustrating a technique for discovering individual characteristic causal relationships. As illustrated in FIG. 9, such a technique first uses an emerging pattern discovery technique to exhaustively find, from a past sample set, a combination of an important factor candidate that has a strong correlation with an objective variable under a specified condition, and a condition at that time. Note that the past sample set is used after being binarized based on a threshold value.

[0004] Thereafter, for each of the found conditions, a causal search technique is used to determine whether the important factor candidate under that condition is accurately an important factor. For example, a case where there is " $x_1 \wedge x_3 \wedge x_4 \rightarrow y$ " ($y = 1$ when $x_1 = x_3 = x_4 = 1$ is true) is assumed. In such a case, one variable selected from a left side is assigned as an "important factor candidate", and the rest is assigned as a "condition". Here, it is assumed that x_4 indicates the "important factor candidate" and the remaining " $x_1 \wedge x_3$ " indicates the "condition". In such a technique, when there is a high correlation between the "important factor candidate" and y on a right side in the past sample set that satisfies the "condition", that "condition" is adopted. The conditions and important factors found in this manner are held in a database (DB). Then, at the time of application, for samples whose causal relationships are desired to be known, the conditions that these samples satisfy are selected from the DB, and the corresponding important factors are presented.

[0005] Non-Patent Document 1: Yusuke Koyanagi, four others, "Developing a Framework for Individual Causal Discovery and its Application to Real Marketing Data", The Japanese Society for Artificial Intelligence 18th Special Interest Group on Business Informatics, March 2021, <URL:[http://sig-bi.jp/doc/18thSIG-BI2021/18thSIG-BI2021 paper13.pdf](http://sig-bi.jp/doc/18thSIG-BI2021/18thSIG-BI2021%20paper13.pdf)>

SUMMARY OF INVENTION

15 TECHNICAL PROBLEM

[0006] However, there is a problem that, when a feature amount included in the past sample set is numerical value data, depending on binarization threshold value setting, a pair of feature amounts in which a high correlation appears under the "condition" may be detected or may not be detected. In other words, when the binarization threshold value setting changes, the correlation for the pair of feature amounts under the "condition" changes. Here, the problem that the correlation for the pair of feature amounts under the "condition" changes when the binarization threshold value setting changes will be described with reference to FIGs. 10A and 10B. FIGs. 10A and 10B are diagrams for describing the problem that the correlation for the pair of feature amounts changes depending on the binarization threshold value setting.

[0007] As illustrated in FIGs. 10A and 10B, in a left diagram, a graph obtained by plotting samples that satisfy the "condition" and samples that do not satisfy the "condition" focusing on a feature amount a and a feature amount b is represented. White circles are the samples that satisfy the "condition", and black circles are the samples that do not satisfy the "condition".

[0008] Under such a situation, in FIG. 10A, a threshold value of the feature amount a when the samples are binarized is set to 3.5, and a threshold value of the feature amount b when the samples are binarized is set to 6.3. Table data of a right diagram is obtained by binarizing the respective samples with these threshold values. In a case where samples that satisfy the "condition" are extracted from binarization, since a correlation between the feature amount a and the feature amount b is high, pairs of the feature amount a and the feature amount b in which a high correlation appears under the "condition" are detected. As a result, this "condition" is adopted as a condition under which a correlation appears.

[0009] On the other hand, in FIG. 10B, the threshold value of the feature amount a when the samples are binarized is set to 5, and the threshold value of the feature amount b when the samples are binarized is set to 4. Table data of a right diagram is obtained by binarizing the respective samples with these threshold values. In a

case where samples that satisfy the "condition" are extracted from binarization, since a correlation between the feature amount a and the feature amount b appears to be low, a pair of the feature amount a and the feature amount b in which a high correlation appears under the "condition" is not detected. As a result, this "condition" is not adopted as a condition under which a correlation appears.

[0010] In this manner, when the feature amount included in the past sample set is the numerical value data, depending on the binarization threshold value setting, a pair of feature amounts in which a high correlation appears under the "condition" may be detected or may not be detected.

[0011] In one aspect, an object of the present invention is to accurately select a condition under which a correlation appears.

SOLUTION TO PROBLEM

[0012] In one aspect, an information processing program causes a computer to execute processing including: generating, from first data in which values of a plurality of attributes included in each sample are accumulated for each sample, second data obtained by binarizing, for each sample, the values of the plurality of attributes included in each sample based on an attribute condition set in advance; enumerating, by using the second data, sets of attribute conditions in which all sample sets indicate true values; computing, for each set of attribute conditions, a correlation between the plurality of attributes in the first data in a sample set associated with each set of attribute conditions; and selecting a set of attribute conditions determined to have a correlation as a condition to be causally searched.

ADVANTAGEOUS EFFECTS OF INVENTION

[0013] According to one embodiment, it is possible to accurately select a condition under which a correlation appears.

BRIEF DESCRIPTION OF DRAWINGS

[0014]

FIG. 1 is a functional block diagram illustrating a configuration of an information processing device according to an embodiment.

FIG. 2 is a diagram illustrating an example of generation processing according to the embodiment.

FIG. 3 is a diagram illustrating an example of enumeration processing according to the embodiment.

FIG. 4A is a diagram illustrating an example of computation processing and selection processing according to the embodiment.

FIG. 4B is a diagram illustrating an example of the computation processing and the selection processing

according to the embodiment.

FIG. 4C is a diagram illustrating an example of the computation processing and the selection processing according to the embodiment.

FIG. 5 is a diagram illustrating an example of a flowchart of information processing according to the embodiment.

FIG. 6 is a diagram illustrating an effect of the information processing according to the embodiment.

FIG. 7 is a diagram illustrating an example of a computer that executes an information processing program.

FIG. 8 is a reference diagram illustrating condition extraction for statistical causal search.

FIG. 9 is a reference diagram illustrating a technique for discovering individual characteristic causal relationships.

FIG. 10A is a diagram for describing a problem that a correlation for a pair of feature amounts changes depending on binarization threshold value setting.

FIG. 10B is a diagram for describing the problem that the correlation for the pair of feature amounts changes depending on the binarization threshold value setting.

DESCRIPTION OF EMBODIMENTS

[0015] Hereinafter, an embodiment of an information processing program, an information processing device, and an information processing method disclosed in the present application will be described in detail with reference to the drawings. Note that the present invention is not limited by the embodiment.

[Embodiment]

[Functional Configuration of Information Processing Device]

[0016] FIG. 1 is a functional block diagram illustrating a configuration of an information processing device according to an embodiment. An information processing device 1 illustrated in FIG. 1 generates, from first data in which values of a plurality of attributes included in each sample are accumulated, second data obtained by binarizing the values of the plurality of attributes included in each sample based on an item. Then, by using the second data, the information processing device 1 enumerates item sets in which all sample sets indicate true values. Then, the information processing device 1 computes a correlation coefficient of the plurality of attribute pairs in the first data in the sample set corresponding to each of the item sets. Then, the information processing device 1 selects an item set predicted to have a correlation as a condition to be causally searched.

[0017] As illustrated in FIG. 1, the information processing device 1 includes a control unit 10 and a storage unit 20. The control unit 10 includes a generation unit 11, an

enumeration unit 12, a computation unit 13, and a selection unit 14.

[0018] The storage unit 20 stores various types of data. The storage unit 20 includes feature amount data 21, binary feature amount data 22, and an item set 23.

[0019] The feature amount data 21 is table data in which values of a plurality of feature amounts (attributes) included in each sample are accumulated. Each row of the feature amount data 21 corresponds to each sample. Each column of the feature amount data 21 corresponds to each feature amount (attribute).

[0020] The binary feature amount data 22 is table data obtained by binarizing values of a plurality of feature amounts (attributes) included in each sample from the feature amount data 21 based on a predetermined item. The binarization is performed by, for example, magnitude comparison between a value of an original feature amount and a constant. Each row of the binary feature amount data 22 corresponds to each sample that is the same as that of the feature amount data 21. Each column of the binary feature amount data 22 corresponds to each binary feature amount. The binary feature amount corresponds to each item for binarizing each feature amount. For example, in a case where the feature amount is a, an item indicating the binary feature amount indicates a condition of magnitude comparison between a and a constant for binarizing the feature amount.

[0021] The item set 23 is groups of item sets enumerated using the binary feature amount data 22, and is groups of item sets corresponding to sample sets that are true with the same binary feature amount (item). The individual item sets included in the item set 23 needs to be frequently appearing item sets. The frequently appearing item set refers to an item set having the number of samples of a certain size or more. The reason why the item set 23 is the frequently appearing item set is that the number of samples having a certain size or more is needed for statistical causal search.

[0022] Furthermore, the individual item sets included in the item set 23 is preferably represented by a saturated item set. The saturated item set refers to a union of groups of item sets having the same sample set. For example, in a case where a sample set {1, 2, 3} has an item set {x3} and also has an item set {x2, x3}, a saturated item set indicating {x2, x3} which is a union of groups of item sets is put as a representative in the item set 23. The reason why the item set 23 is the saturated item set is that it is useless to associate a plurality of item sets with the same sample set. Note that the item set 23 is generated by the enumeration unit 12.

[0023] The generation unit 11 generates the binary feature amount data 22 from the feature amount data 21. For example, the generation unit 11 generates, from the feature amount data 21, the binary feature amount data 22 obtained by binarizing values of the respective feature amounts into true and false based on binary feature amounts (items). Then, the generation unit 11 stores the generated binary feature amount data 22 in the

storage unit 20.

[0024] The enumeration unit 12 enumerate, by using the binary feature amount data 22, frequently appearing saturated item sets in which all sample sets indicate true values. For example, by using the binary feature amount data 22, the enumeration unit 12 extracts correspondence information in which all sample sets indicate true values are associated with each item set. Each item set is a frequently appearing item set having the number of samples of a certain size or more. Additionally, each item set is a saturated item set that is a union of groups of item sets having the same sample set. That is, each item set is the frequently appearing saturated item set. Then, the enumeration unit 12 stores a group of frequently appearing saturated item sets in the storage unit 20 as the item set 23.

[0025] Furthermore, the enumeration unit 12 generates a directed acyclic graph having an empty set as a starting point for the enumerated frequently appearing saturated item sets. In the directed acyclic graph, items are added from upstream toward downstream. Note that the frequently appearing saturated item sets may be enumerated using, for example, an algorithm of enumeration of the frequently appearing saturated item sets of "Takeaki Uno, Tatsuya Asai, Yuzo Uchida, Hiroki Arimura, "An Efficient Algorithm for Enumerating Closed Patterns in Transaction Databases", Discovery Science 2004, LNAI 3245, pp.16-31".

[0026] The computation unit 13 computes a correlation between a plurality of feature amounts in the feature amount data 21 in a sample set associated with each frequently appearing saturated item set. For example, the computation unit 13 selects a frequently appearing saturated item set in order from an empty set of a directed acyclic graph toward downstream. The computation unit 13 computes a correlation coefficient of each feature amount pair in the feature amount data 21 in a sample set that satisfies conditions of all items included in the selected frequently appearing saturated item set. As an example, the computation unit 13 extracts, from the feature amount data 21, a value of a feature amount of each sample set associated with the selected frequently appearing saturated item set. Then, the computation unit 13 computes a correlation coefficient of each feature amount pair using the extracted value of the feature amount.

[0027] The selection unit 14 selects a frequently appearing saturated item set predicted to have a correlation as a condition to be causally searched. For example, in a case where there are a certain number or more of feature amount pairs having a correlation coefficient equal to or larger than a threshold value set for the selected frequently appearing saturated item set, the selection unit 14 determines that there is a correlation and selects the selected frequently appearing saturated item set as a condition. The certain number may be "1" or "2", and it is sufficient the certain number is set in advance.

[0028] Note that the threshold value is set for each

feature amount pair. Then, the threshold value is updated such that there is a feature amount pair whose correlation increases by a certain amount or more due to addition of an item from upstream to downstream. Furthermore, the threshold value may be updated such that a set of feature amount pairs whose correlation increases by a certain amount or more due to addition of an item is different from others.

[Example of Generation Processing]

[0029] Here, generation processing performed by the generation unit 11 according to the embodiment will be described with reference to FIG. 2. FIG. 2 is a diagram illustrating an example of the generation processing according to the embodiment. As illustrated in FIG. 2, in a left diagram, feature amount data Y is represented as the feature amount data 21. In the feature amount data Y, a, b, c,... are set as feature amounts, and value of the feature amounts are set for each sample.

[0030] The generation unit 11 generates the binary feature amount data 22 from such feature amount data Y. In other words, the generation unit 11 generates, from the feature amount data Y, the binary feature amount data 22 obtained by binarizing the values of the respective feature amounts into true and false based on binary feature amounts (items) set in advance. In a right diagram, binary feature amount data X is represented as the binary feature amount data 22. In the binary feature amount data X, x_1 , x_2 , x_3 , ... are set as the binary feature amounts (items), and "1" and "0" are set for the respective samples. The x_1 as the item indicates a condition that the feature amount a is smaller than 5. The x_2 as the item indicates a condition that the feature amount a is larger than 2. The x_3 indicates a condition that the feature amount c is 0. Then, for each of the x_1 , x_2 , and x_3 , "1" is set in a case where the condition is satisfied, and "0" is set in a case where the condition is not satisfied.

[0031] As an example, in a case where the sample is "1", for the item x_1 ($a < 5$), "1" is set because the feature amount a is "1.3" and the feature amount a is smaller than "5". For the item x_2 ($a \geq 2$), "0" is set because the feature amount a is "1.3" and the feature amount a is smaller than "2". For the item x_3 ($c = 0$), "0" is set because the feature amount c is "1" and the feature amount c is not "0".

[Example of Enumeration Processing]

[0032] Next, enumeration processing performed by the enumeration unit 12 according to the embodiment will be described with reference to FIG. 3. FIG. 3 is a diagram illustrating an example of the enumeration processing according to the embodiment. As illustrated in FIG. 3, in a left diagram, the binary feature amount data X is represented as the binary feature amount data 22.

[0033] By using such binary feature amount data X, the enumeration unit 12 extracts item sets in which all sample sets indicate true values. Here, it is assumed that the

extracted item set has the number of three or more samples. The enumeration unit 12 extracts the following item sets. An empty set $\{\}$ in a case where the sample set is $\{1, 2, 3, 4, 5, 6, 7\}$ is extracted. $\{x_1\}$ in a case where the sample set is $\{1, 2, 3, 4\}$ is extracted. $\{x_2\}$ in a case where the sample set is $\{2, 3, 4, 5, 6\}$ is extracted. $\{x_2, x_3\}$ in a case where the sample set is $\{2, 3, 4, 5\}$ is extracted. $\{x_1, x_2\}$ in a case where the sample set is $\{2, 3, 4\}$ is extracted. $\{x_1, x_3\}$ in a case where the sample set is $\{2, 3, 4\}$ is extracted. $\{x_1, x_2, x_3\}$ in a case where the sample set is $\{2, 3, 4\}$ is extracted.

[0034] Then, the enumeration unit 12 enumerates frequently appearing saturated item sets from the extracted item sets. Here, since the item sets $\{x_1, x_2\}$, $\{x_1, x_3\}$, and $\{x_1, x_2, x_3\}$ indicated by a reference sign i1 have the same sample set $\{2, 3, 4\}$, the $\{x_1, x_2, x_3\}$ is exemplified as the frequently appearing saturated item set. Furthermore, since the item sets $\{x_3\}$ and $\{x_2, x_3\}$ indicated by a reference sign i2 have the same sample set $\{2, 3, 4, 5\}$, the $\{x_2, x_3\}$ is exemplified as the frequently appearing saturated item set. Note that the empty set $\{\}$, the $\{x_1\}$, and the $\{x_2\}$ are similarly exemplified as the frequently appearing saturated item sets.

[0035] Then, the enumeration unit 12 generates a directed acyclic graph having the empty set as a starting point for the enumerated frequently appearing saturated item sets. In a right diagram, the directed acyclic graph for these frequently appearing saturated item sets is represented.

[Example of Computation Processing and Selection Processing]

[0036] Next, computation processing and selection processing performed by the computation unit 13 and the selection unit 14 according to the embodiment will be described with reference to FIGs. 4A to 4C. FIGs. 4A to 4C are diagrams illustrating an example of the computation processing and the selection processing according to the embodiment. Note that, in FIGs. 4A to 4C, a correlation coefficient of each feature amount pair in the feature amount data 21 for feature amounts a, b, c, and d is computed.

[0037] As illustrated in FIG. 4A, in an upper left diagram, initial values of a positive correlation threshold value and a negative correlation threshold value of each feature amount pair for the feature amounts a, b, c, and d are set. Here, an initial value of a positive correlation threshold value θ indicates 0.4. An initial value of a negative correlation threshold value θ indicates -0.4.

[0038] Under such a situation, the computation unit 13 selects an item set indicating a starting point from a directed acyclic graph for frequently appearing saturated item sets. Then, the computation unit 13 computes a correlation coefficient of each feature amount pair in the feature amount data 21 in a sample set that satisfies a condition of an empty set $\{\}$ that is the selected item set. The sample set that satisfies the condition of the empty

set $\{\}$ mentioned here means the entire samples. Here, the correlation coefficient of each feature amount pair in a case where the item set is the empty set $\{\}$ is represented in a right diagram. In the correlation coefficient of each feature amount pair in the right diagram, a positive correlation coefficient of a pair of the feature amounts a and b is 0.7, which is larger than the threshold value 0.4. Furthermore, a negative correlation coefficient of a pair of the feature amounts d and b is -0.5, which is smaller than the threshold value -0.4.

[0039] Thus, since there is the feature amount pair having the correlation coefficient equal to or larger than the set threshold value for the positive correlation, the selection unit 14 adopts (selects) the empty set $\{\}$ that is the selected item set as a condition. Furthermore, since there is also the feature amount pair having the correlation coefficient equal to or smaller than the set threshold value from the negative correlation, the selection unit 14 may adopt (select) the empty set $\{\}$ that is the selected item set as a condition.

[0040] Then, the selection unit 14 updates the threshold value used downstream of the empty set $\{\}$ that is the item set. In other words, the selection unit 14 updates the threshold value in order to find, as a condition, an item set better than the item set already adopted (selected) as the condition. The updated threshold value is represented in a lower left diagram. Here, a feature amount pair having a high correlation in a case where the item set is the empty set $\{\}$ is the pair of feature amounts a and b. Therefore, the positive correlation threshold value of the pair of feature amounts a and b is updated so as to increase by a certain amount δ ($= 0.2$). In other words, the positive correlation threshold value of the pair of feature amounts a and b is updated to 0.9. That is, such a threshold value is updated such that there is a feature amount pair whose correlation increases by the certain amount δ ($= 0.2$) or more due to addition of an item from the empty set $\{\}$ to downstream. Furthermore, a feature amount pair having a high correlation in a case where the item set is the empty set $\{\}$ is the pair of feature amounts d and b. Therefore, the negative correlation threshold value of the pair of feature amounts d and b is updated so as to decrease by the certain amount δ ($= 0.2$). In other words, the negative correlation threshold value of the pair of feature amounts d and b is updated to -0.7. That is, such a threshold value is updated such that there is a feature amount pair whose correlation increases by the certain amount δ ($= 0.2$) or more due to addition of an item from the empty set $\{\}$ to downstream. Note that, here, the positive correlation coefficient of a pair of the feature amounts c and d is 0.3, which is lower than the threshold value 0.4, but a positive correlation threshold value of the pair of feature amounts c and d is updated so as to increase by the certain amount δ ($= 0.2$) in order to find a better item set. In other words, the positive correlation threshold value of the pair of feature amounts c and d is updated to 0.5. That is, such a threshold value is updated such that a set of feature amount pairs whose correlation increases by the certain

amount or more due to addition of an item from the empty set $\{\}$ to downstream is different from others.

[0041] As illustrated in FIG. 4B, in a left diagram, the threshold values updated from the empty set $\{\}$ to downstream are set. Here, the positive correlation threshold value of the pair of feature amounts a and b is updated to 0.9. The negative correlation threshold value of the pair of feature amounts d and b is updated to -0.7. Furthermore, the positive correlation threshold value of the pair of feature amounts c and d is updated to 0.5.

[0042] Under such a situation, the computation unit 13 selects the next item set from the directed acyclic graph for the frequently appearing saturated item sets. Here, it is assumed that $\{x_1\}$ is selected as an item set downstream of the empty set $\{\}$. Then, the computation unit 13 computes a correlation coefficient of each feature amount pair in the feature amount data 21 in a sample set that satisfies a condition of the selected item set $\{x_1\}$. For example, the correlation coefficient of each feature amount pair in a case where the item set is $\{x_1\}$ is represented in a right diagram. In the correlation coefficient of each feature amount pair in the right diagram, the positive correlation coefficients of all the feature amount pairs are lower than the threshold values updated from the empty set $\{\}$ to downstream. The negative correlation coefficients of all the feature amount pairs are higher than the threshold values updated from the empty set $\{\}$ to downstream.

[0043] Thus, the selection unit 14 does not adopt (select) the item set $\{x_1\}$ as a condition. That is, the sample set that satisfies the item set $\{x_1\}$ is a subset of a sample set that satisfies the empty set $\{\}$, which is a similar sample set, and there is no feature amount pair in which a correlation increases by the certain amount δ or more. Therefore, the item set $\{x_1\}$ is not adopted (selected) as a condition.

[0044] As illustrated in FIG. 4C, in an upper left diagram, the threshold values updated from the empty set $\{\}$ to downstream are set. The updated threshold values are the same as those indicated in the left diagram of FIG. 4B. In other words, the positive correlation threshold value of the pair of feature amounts a and b is updated to 0.9. The negative correlation threshold value of the pair of feature amounts d and b is updated to -0.7. Furthermore, the positive correlation threshold value of the pair of feature amounts c and d is updated to 0.5.

[0045] Under such a situation, the computation unit 13 selects the next item set from the directed acyclic graph for the frequently appearing saturated item sets. Here, it is assumed that $\{x_2\}$ is selected as an item set downstream of the empty set $\{\}$. Then, the computation unit 13 computes a correlation coefficient of each feature amount pair in the feature amount data 21 in a sample set that satisfies a condition of the selected item set $\{x_2\}$. For example, the correlation coefficient of each feature amount pair in a case where the item set is $\{x_2\}$ is represented in a right diagram. In the correlation coefficient of each feature amount pair in the right diagram, a

positive correlation coefficient of a pair of the feature amounts a and d is 0.7, which is larger than the threshold value 0.4.

[0046] Thus, since there is the feature amount pair having the correlation coefficient equal to or larger than the set threshold value for the positive correlation, the selection unit 14 adopts (selects) the selected item set $\{x_2\}$ as a condition.

[0047] Then, the selection unit 14 updates the threshold value used downstream of the item set $\{x_2\}$. In other words, the selection unit 14 updates the threshold value in order to find, as a condition, an item set better than the item set already adopted (selected) as the condition. The updated threshold value is represented in a lower left diagram. Here, a feature amount pair having a high correlation in a case where the item set is $\{x_2\}$ is the pair of feature amounts a and d . Therefore, the positive correlation threshold value of the pair of feature amounts a and d is updated so as to increase by the certain amount δ ($= 0.2$). In other words, the positive correlation threshold value of the pair of feature amounts a and d is updated to 0.9. That is, such a threshold value is updated such that there is a feature amount pair whose correlation increases by the certain amount or more due to addition of an item from upstream to downstream. Additionally, threshold values of the other feature amount pairs indicate the threshold values updated by the empty set $\{\}$ indicated in the upper left diagram. Note that, here, the negative correlation coefficient of the pair of feature amounts d and b is -0.6, which is lower than the threshold value -0.7, but the negative correlation threshold value of the pair of feature amounts d and b is updated so as to decrease by the certain amount δ ($= 0.2$) in order to find a better item set. In other words, the negative correlation threshold value of the pair of feature amounts d and b is updated to -0.8. That is, such a threshold value is updated such that a set of feature amount pairs whose correlation increases by the certain amount or more due to addition of an item from upstream to downstream is different from others.

[Flowchart of Information Processing]

[0048] FIG. 5 is a diagram illustrating an example of a flowchart of information processing according to the embodiment.

[0049] As illustrated in FIG. 5, the generation unit 11 binarizes the feature amount data Y to generate the binary feature amount data X (step S11).

[0050] Then, the enumeration unit 12 enumerates frequently appearing saturated item sets in the binary feature amount data X , and generates a directed acyclic graph having an empty set $\{\}$ as a starting point for the enumerated frequently appearing saturated item sets (step S12).

[0051] Then, the computation unit 13 selects a frequently appearing saturated item set I by giving priority to depth or giving priority to width from the empty set $\{\}$

(step S13). The computation unit 13 computes a correlation coefficient of each feature amount pair in the feature amount data Y in a sample set that satisfies a condition of the frequently appearing saturated item set I (step S14).

[0052] Then, when there is a feature amount pair having a correlation coefficient equal to or larger than a threshold value set upstream of the frequently appearing saturated item set I , the selection unit 14 selects and outputs the selected frequently appearing saturated item set I as a condition (step S15). Note that, in a case where the frequently appearing saturated item set I is the empty set at the starting point, the threshold value used in the frequently appearing saturated item set I is assumed to be set in advance.

[0053] Then, the selection unit 14 sets a threshold value to be used downstream of the selected frequently appearing saturated item set I (step S16). For example, the selection unit 14 updates a threshold value of a feature amount pair having a correlation equal to or greater than the threshold value, so as to increase the threshold value by a certain amount or more. Furthermore, the selection unit 14 may update a threshold value of a feature amount pair different from the feature amount pair having the correlation equal to or greater than the threshold value, so as to increase the threshold value by a certain amount or more.

[0054] Then, the computation unit 13 proceeds to step S13 so as to repeat until there is no unselected frequently appearing saturated item set (step S17). Then, when there is no unselected frequently appearing saturated item set, the computation unit 13 ends the information processing.

[0055] In this manner, the information processing according to the embodiment may accurately select a condition under which a correlation appears.

[0056] FIG. 6 is a diagram illustrating an effect of the information processing according to the embodiment. In a left diagram, a graph obtained by plotting samples that satisfy the "condition" and samples that do not satisfy the "condition" focusing on a feature amount a and a feature amount b is represented. White circles are the samples that satisfy the "condition", and black circles are the samples that do not satisfy the "condition". In the information processing according to the embodiment, a pair of feature amounts is evaluated not from information binarized by threshold value setting but from information indicating values of original feature amounts. Therefore, a pair of feature amounts in which a high correlation appears under the "condition" is not overlooked. Therefore, the information processing according to the embodiment may accurately select the "condition" under which a correlation appears based on a correlation coefficient between feature amounts computed from values of original feature amounts.

[Effects of Embodiment]

[0057] According to the embodiment described above,

the information processing device 1 generates, from the feature amount data 21 in which values of a plurality of feature amounts included in each sample are accumulated for each sample, the binary feature amount data 22 obtained by binarizing, for each sample, the values of the plurality of feature amounts included in each sample based on an item set in advance. By using the binary feature amount data 22, the information processing device 1 enumerates item sets in which all sample sets indicate true values. Then, for each item set, the information processing device 1 computes a correlation between a plurality of feature amounts in the feature amount data 21 in a sample set associated with each item set. Then, the information processing device 1 selects an item set determined to have a correlation as a condition to be causally searched. As a result, the information processing device 1 may accurately select a condition under which a correlation appears.

[0058] Furthermore, according to the embodiment described above, the information processing device 1 enumerates item sets indicating a union of groups of item sets having the same sample set. As a result, the information processing device 1 may suppress extraction of a plurality of item sets having the same sample set.

[0059] Furthermore, according to the embodiment described above, the information processing device 1 further enumerates item sets for a sample set in which the number of samples included in the sample set is equal to or larger than a predetermined number. As a result, the information processing device 1 may increase accuracy of selecting a condition as the number of samples included in the sample set is equal to or larger than the predetermined number.

[0060] Furthermore, according to the embodiment described above, the information processing device 1 selects an item set in order, and computes a correlation coefficient of a plurality of pairs of feature amounts in the feature amount data 21 in a sample set associated with the selected item set. In a case where there is a pair of feature amounts having a correlation coefficient equal to or larger than a predetermined threshold value, the information processing device 1 selects the selected item set as a condition. As a result, the information processing device 1 may avoid overlooking a pair of feature amounts in which a high correlation appears under the condition by using, for evaluation of the correlation coefficient, values of original feature amounts in the sample set.

[0061] Furthermore, according to the embodiment described above, the information processing device 1 generates a directed acyclic graph having an item set as an empty set as a starting point by using the enumerated item sets. The information processing device 1 selects an item set in order by giving priority to depth or giving priority to width from the empty set included in the directed acyclic graph, and computes a correlation coefficient related to the selected item set. Then, in a case where there is a pair of feature amounts having a correlation coefficient equal to or larger than a threshold value set at

a higher level of the selected item set, the information processing device 1 selects the selected item set as a condition. As a result, the information processing device 1 may select a better item set as a condition by using the threshold value set at the higher level.

[0062] Note that each illustrated component of the information processing device 1 does not necessarily have to be physically configured as illustrated in the drawings. In other words, specific forms of distribution and integration of the information processing device 1 are not limited to the illustrated ones, and the whole or a part of the information processing device 1 may be configured by being functionally or physically distributed and integrated in optional units according to various loads, use situations, or the like. For example, the computation unit 13 and the selection unit 14 may be integrated. Furthermore, the storage unit 20 may be coupled through a network as an external device of the information processing device 1.

[0063] Furthermore, various types of processing described in the embodiment described above may be implemented by a computer such as a personal computer or a workstation executing programs prepared in advance. Thus, in the following, an example of a computer that executes an information processing program that implements functions similar to the functions of the information processing device 1 illustrated in FIG. 1 will be described. Here, the information processing program that implements the functions similar to the functions of the information processing device 1 will be described as an example. FIG. 7 is a diagram illustrating an example of the computer that executes the information processing program.

[0064] As illustrated in FIG. 7, a computer 200 includes a CPU 203 that executes various types of computation processing, an input device 215 that accepts data input from a user, and a display control unit 207 that controls a display device 209. Furthermore, the computer 200 includes a drive device 213 that reads a program and the like from a storage medium, and a communication control unit 217 that exchanges data with another computer via a network. Furthermore, the computer 200 includes a memory 201 that temporarily stores various types of information, and a hard disk drive (HDD) 205. Additionally, the memory 201, the CPU 203, the HDD 205, the display control unit 207, the drive device 213, the input device 215, and the communication control unit 217 are coupled by a bus 219.

[0065] The drive device 213 is, for example, a device for a removable disk 210. The HDD 205 stores an information processing program 205a and information processing related information 205b.

[0066] The CPU 203 reads the information processing program 205a, loads the information processing program 205a into the memory 201, and executes the information processing program 205a as a process. Such a process corresponds to each functional unit of the information processing device 1. The information processing related

information 205b corresponds to the feature amount data 21, the binary feature amount data 22, and the item set 23. Additionally, for example, the removable disk 210 stores each piece of information such as the information processing program 205a.

[0067] Note that the information processing program 205a does not necessarily have to be stored in the HDD 205 from the beginning. For example, the program is stored in a "portable physical medium" to be inserted into the computer 200, such as a flexible disk (FD), a compact disk read only memory (CD-ROM), a digital versatile disk (DVD), a magneto-optical disk, or an integrated circuit (IC) card. Then, the computer 200 may read the information processing program 205a from these media to execute the information processing program 205a.

REFERENCE SIGNS LIST

[0068]

1 Information processing device

10 Control unit

11 Generation unit

12 Enumeration unit

13 Computation unit

14 Selection unit

20 Storage unit

21 Feature amount data

22 Binary feature amount data

23 Item set

Claims

1. An information processing program for causing a computer to execute processing comprising:

generating, from first data in which values of a plurality of attributes included in each sample are accumulated for each sample, second data obtained by binarizing, for each sample, the values of the plurality of attributes included in each sample based on an attribute condition set in advance;

enumerating, by using the second data, sets of attribute conditions in which all sample sets indicate true values;

computing, for each set of attribute conditions, a correlation between the plurality of attributes in the first data in a sample set associated with each set of attribute conditions; and

selecting a set of attribute conditions determined to have a correlation as a condition to be causally searched.

2. The information processing program according to claim 1, wherein,
in the processing of enumerating the sets of attribute

conditions, the sets of attribute conditions that indicate a union of groups of sets of attribute conditions that have the same sample set are enumerated.

3. The information processing program according to claim 2, wherein,
in the processing of enumerating the sets of attribute conditions, the sets of attribute conditions are further enumerated for a sample set in which the number of samples included in the sample set is equal to or larger than a predetermined number.

4. The information processing program according to claim 1, wherein,

in the processing of computing the correlation, the set of attribute conditions is selected in order, and a correlation coefficient of a plurality of pairs of attributes in the first data in a sample set associated with the selected set of attribute conditions is computed, and
in the processing of selecting the set of attribute conditions, in a case where there is a pair of attributes that has a correlation coefficient equal to or larger than a predetermined threshold value, the selected set of attribute conditions is selected as the condition.

5. The information processing program according to claim 4, wherein,

in the processing of enumerating the sets of the attribute conditions, a directed acyclic graph that has the set of attribute conditions as an empty set as a starting point is generated by using the enumerated sets of attribute conditions,
in the processing of computing the correlation, the set of attribute conditions is selected in order by giving priority to depth or giving priority to width from the empty set included in the directed acyclic graph, and a correlation coefficient related to the selected set of attribute conditions is computed, and
in the processing of selecting the set of attribute conditions, in a case where there is a pair of attributes that has a correlation coefficient equal to or larger than a threshold value set at a higher level of the selected set of attribute conditions, the selected set of attribute conditions is selected as the condition.

6. An information processing device comprising:

a generation unit configured to generate, from first data in which values of a plurality of attributes included in each sample are accumulated for each sample, second data obtained by binarizing, for each sample, the values of the

plurality of attributes included in each sample
 based on an attribute condition set in advance;
 an enumeration unit configured to enumerate,
 by using the second data, sets of attribute con-
 ditions in which all sample sets indicate true 5
 values;
 a computation unit configured to compute, for
 each set of attribute conditions, a correlation
 between the plurality of attributes in the first data
 in a sample set associated with each set of 10
 attribute conditions; and
 a selection unit configured to select a set of
 attribute conditions determined to have a corre-
 lation as a condition to be causally searched.

15

7. An information processing method for causing a
 computer to execute processing comprising:

generating, from first data in which values of a
 plurality of attributes included in each sample 20
 are accumulated for each sample, second data
 obtained by binarizing, for each sample, the
 values of the plurality of attributes included in
 each sample based on an attribute condition set
 in advance; 25
 enumerating, by using the second data, sets of
 attribute conditions in which all sample sets
 indicate true values;
 computing, for each set of attribute conditions, a
 correlation between the plurality of attributes in 30
 the first data in a sample set associated with
 each set of attribute conditions; and
 selecting a set of attribute conditions deter-
 mined to have a correlation as a condition to
 be causally searched. 35

40

45

50

55

FIG. 1

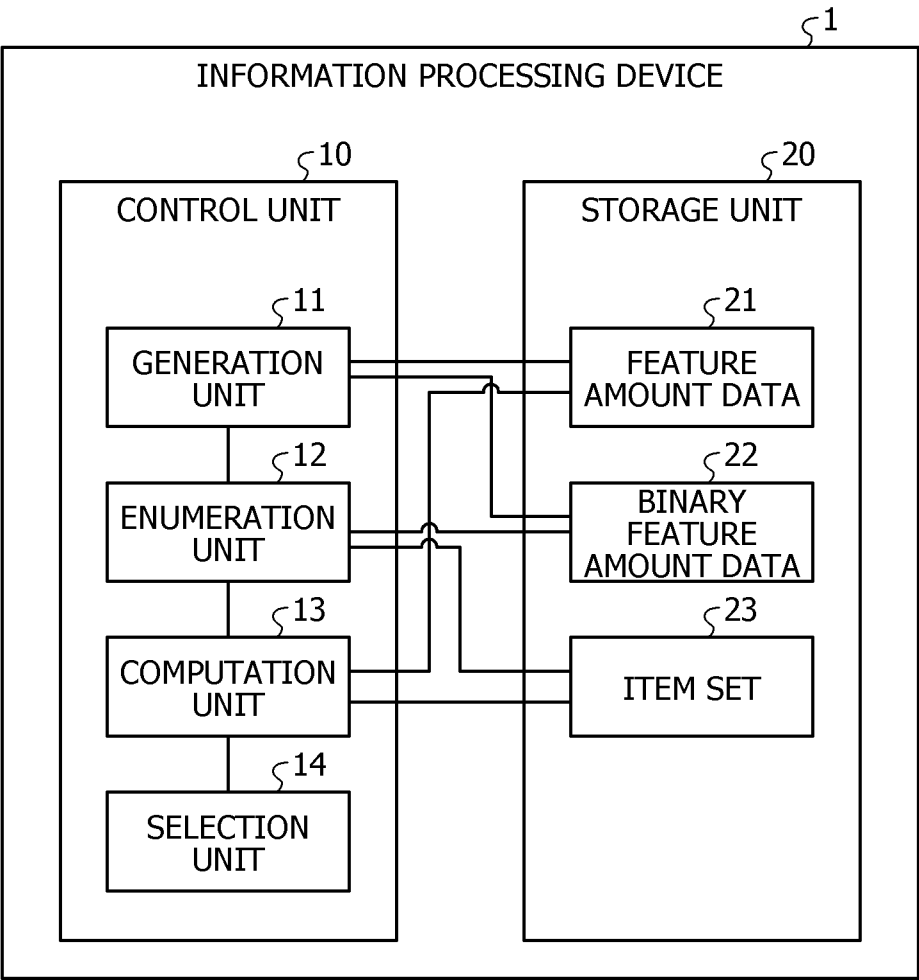


FIG. 2

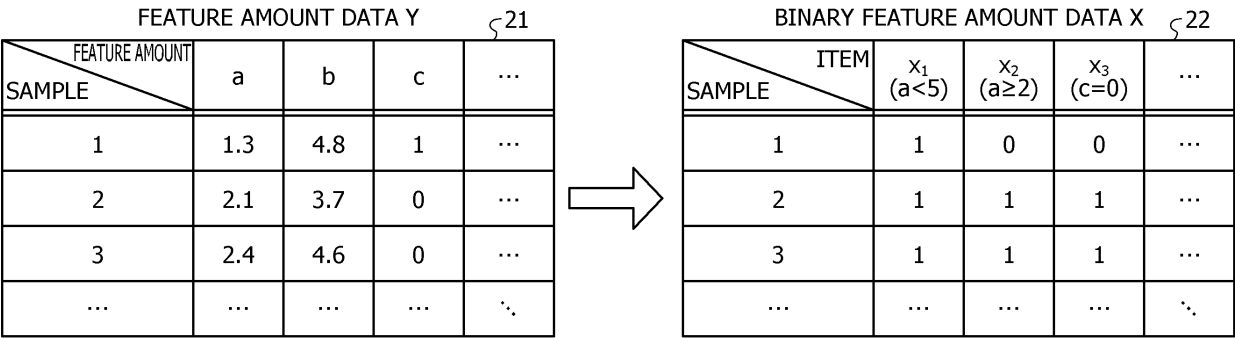


FIG. 3

THREE OR MORE
FREQUENTLY APPEAR

§ 22

BINARY FEATURE AMOUNT DATA X				
SAMPLE	ITEM	x ₁	x ₂	x ₃
1		1	0	0
2		1	1	1
3		1	1	1
4		1	1	1
5		0	1	1
6		0	1	0
7		0	0	1

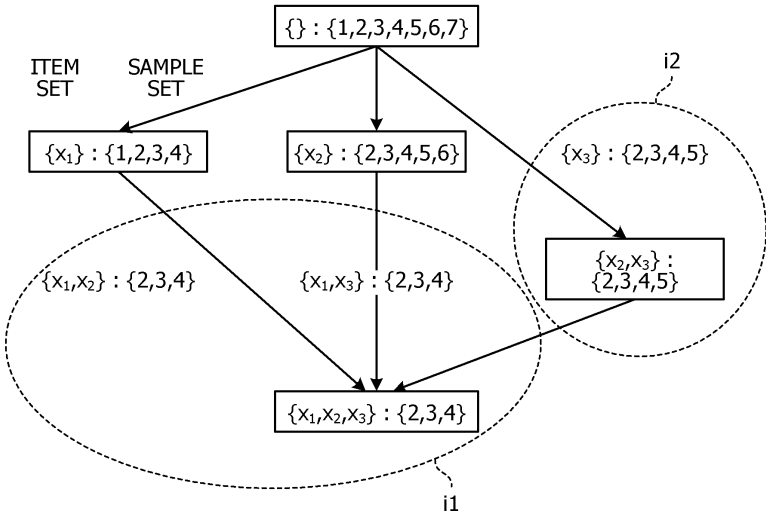


FIG. 4A

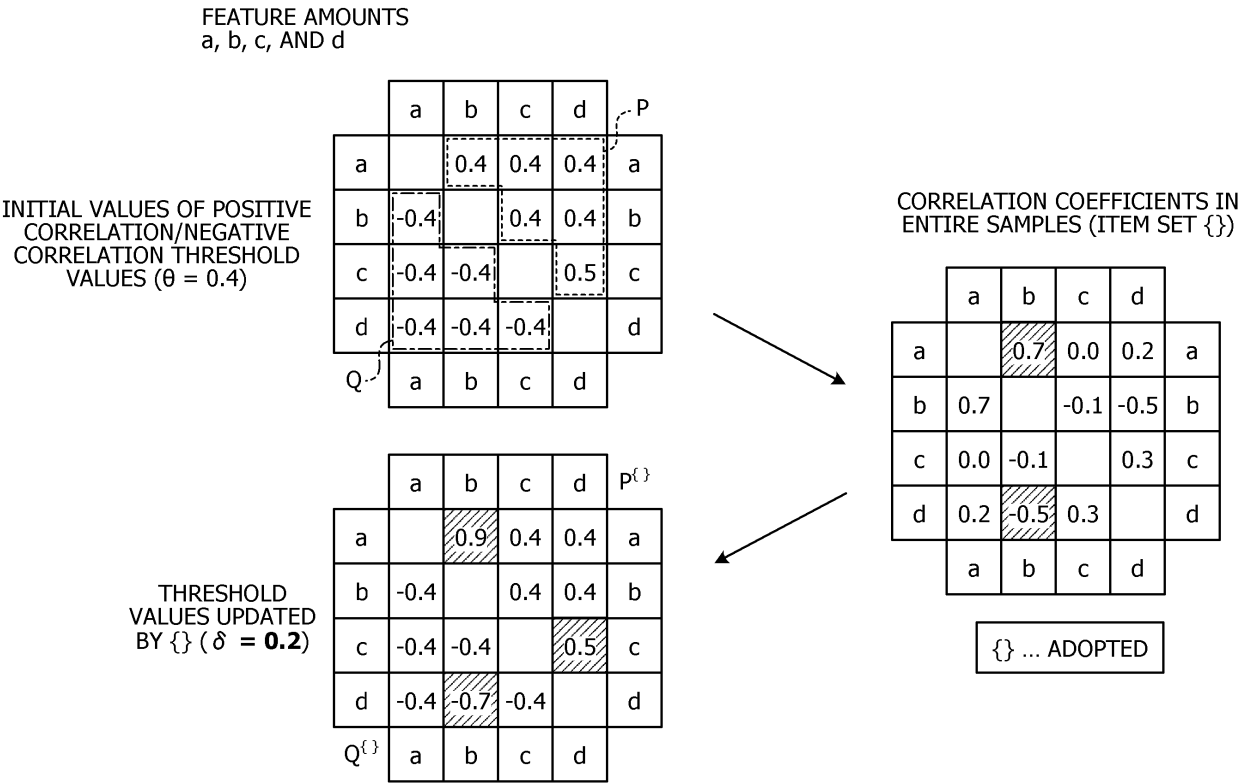


FIG. 4B

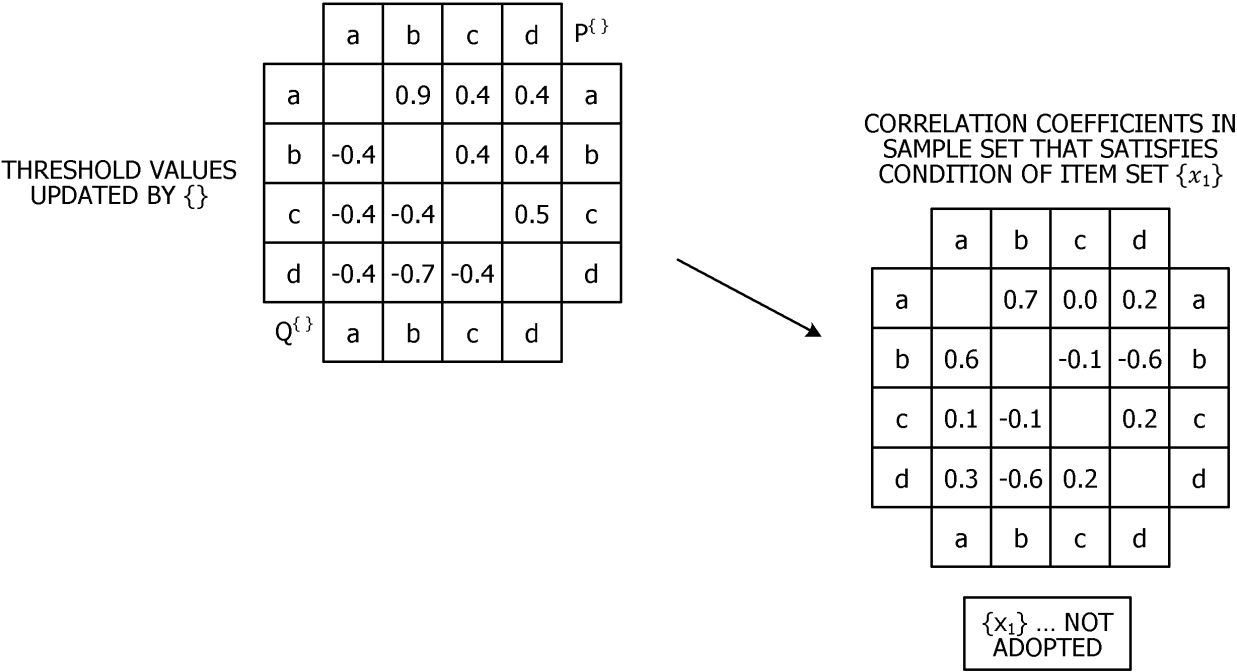


FIG. 4C

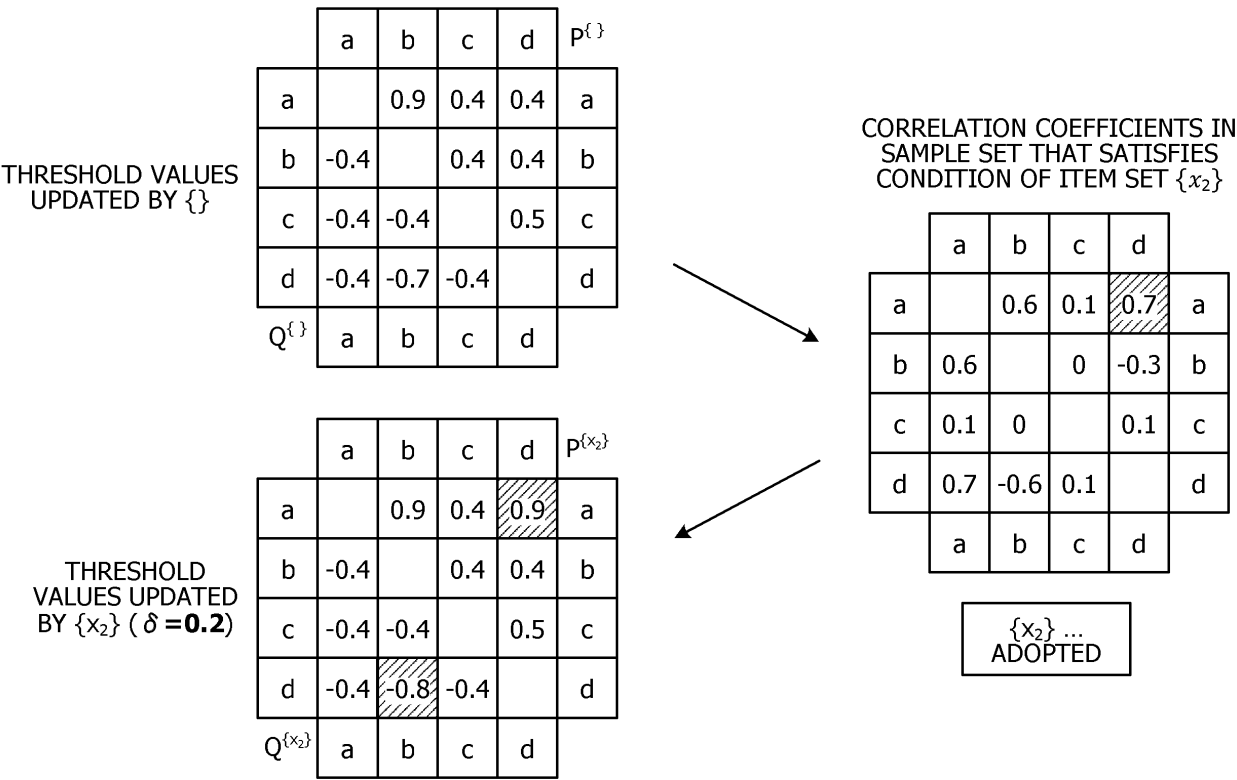


FIG. 5

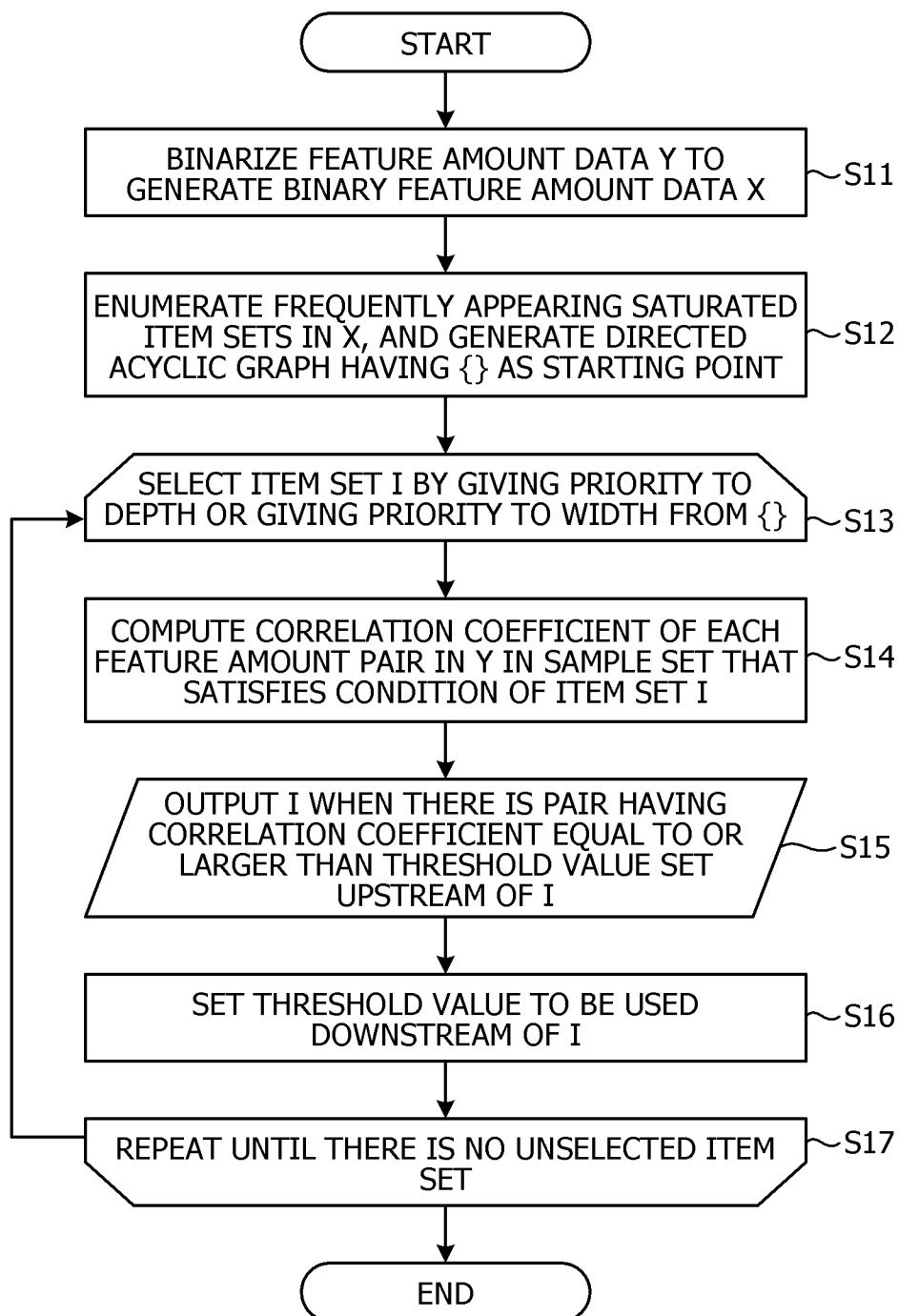


FIG. 6

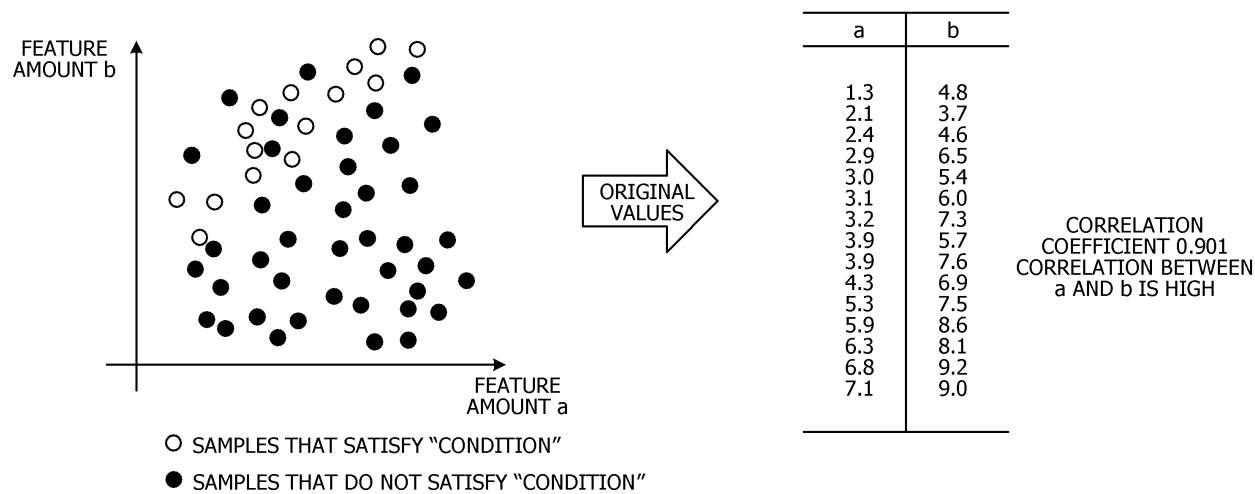


FIG. 7

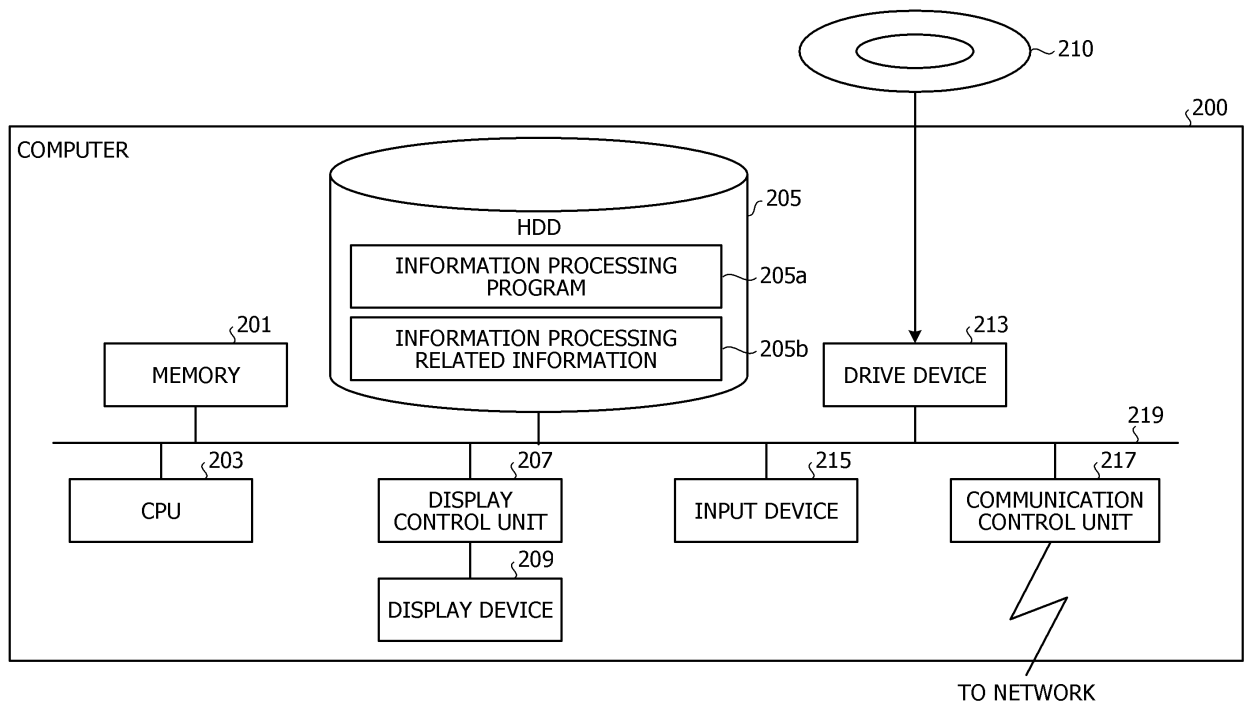


FIG. 8

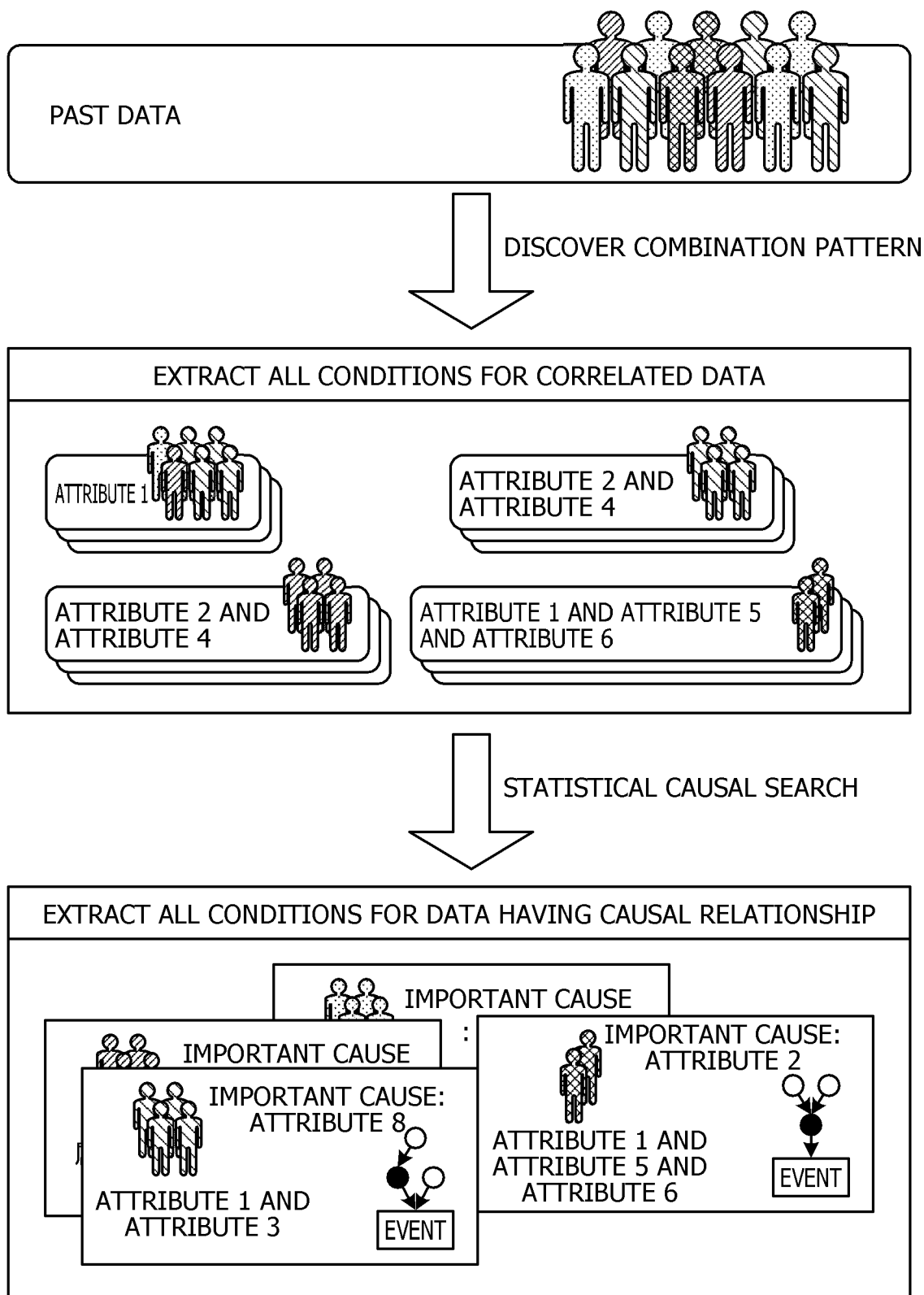


FIG. 9

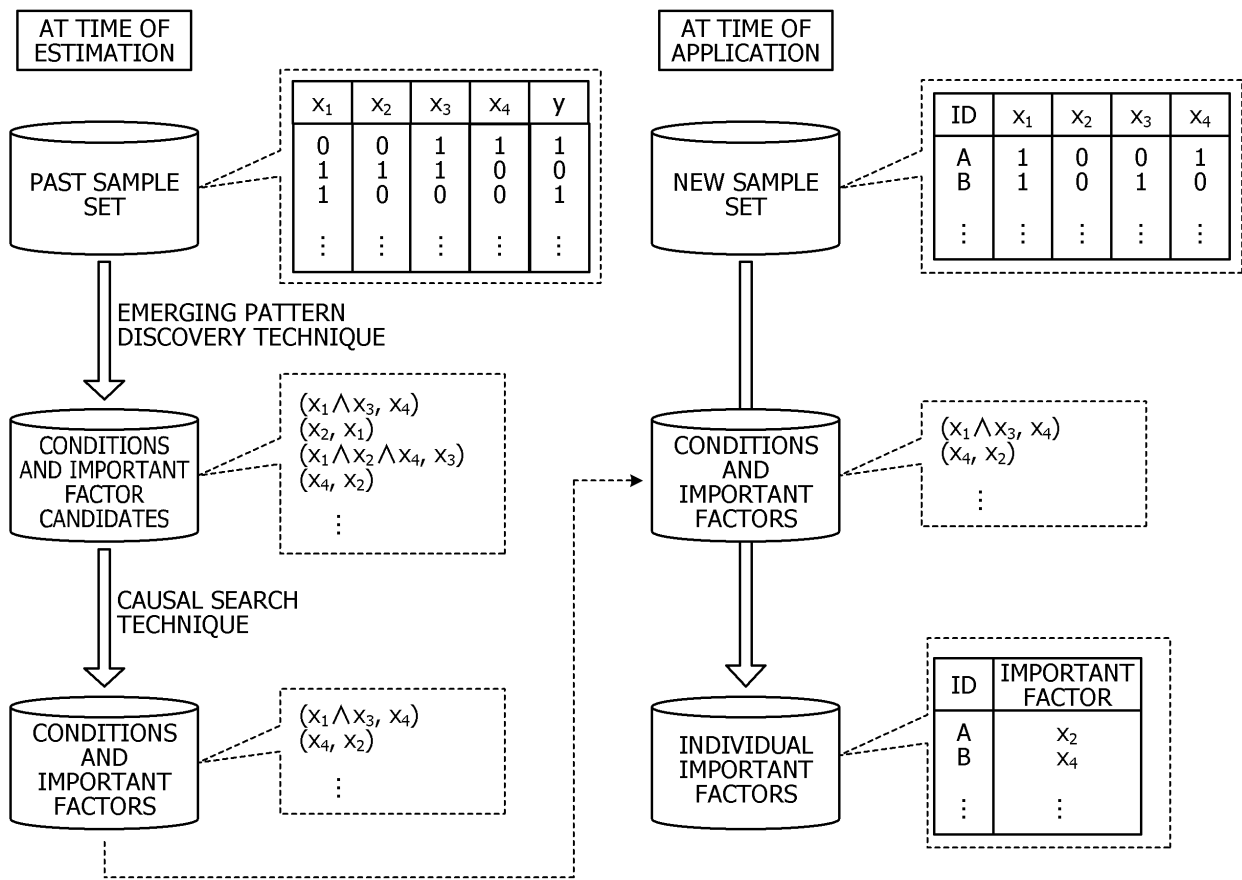


FIG. 10A

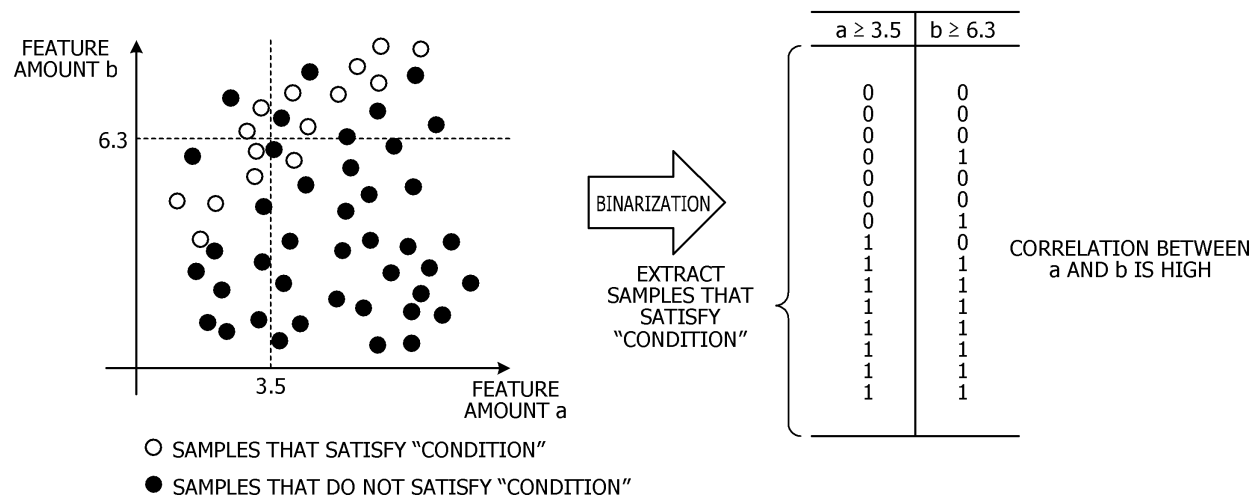
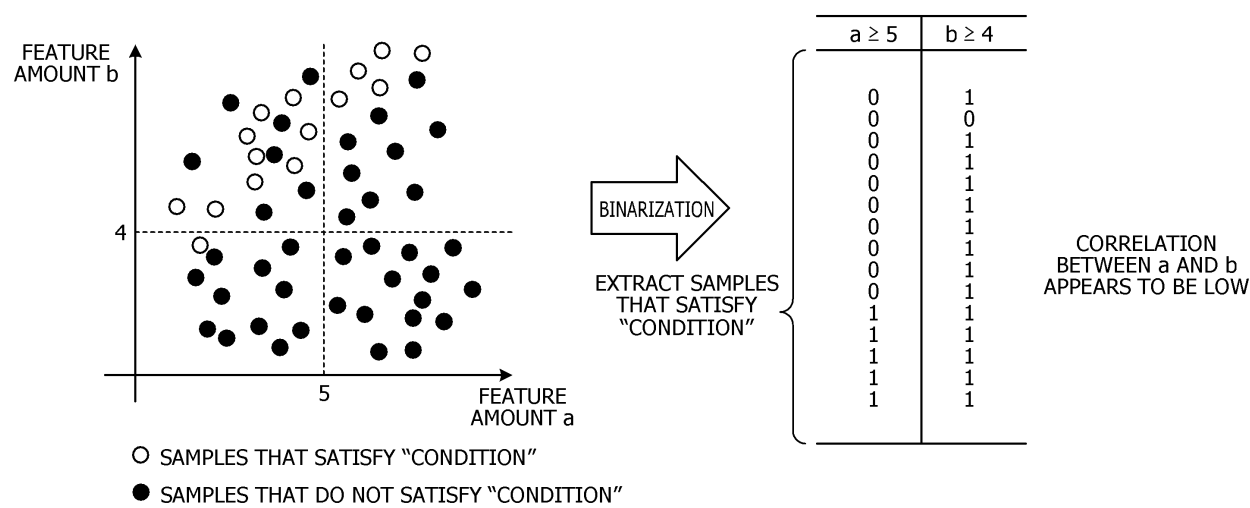


FIG. 10B



INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2022/005497

A. CLASSIFICATION OF SUBJECT MATTER

G06N 5/02(2006.01)i

FI: G06N5/02 150

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

G06N5/02

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Published examined utility model applications of Japan 1922-1996
 Published unexamined utility model applications of Japan 1971-2022
 Registered utility model specifications of Japan 1996-2022
 Published registered utility model applications of Japan 1994-2022

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	小柳 佑介 ほか, 個々の特徴的な因果関係を発見する技術の開発とマーケティングデータへの適用, S I G - B I 人工知能学会: 経営課題に A I を! ビジネス・インフォマティクス研究会 第 1 8 回研究会, 05 March 2021 entire text, all drawings, (KOYANAGI, Yusuke et al. Developing a Framework for Individual Causal Discovery and its Application to Real Marketing Data.), non-official translation (SIG-BI JSAI: Application AI to management issue! Special Interest Group on Business Informatics. The 18th Research Meeting.)	1-7
A	JP 2001-265596 A (MITSUBISHI ELECTRIC CORP) 28 September 2001 (2001-09-28) entire text, all drawings	1-7
A	US 2019/0384800 A1 (ELASSAAD, Shauki) 19 December 2019 (2019-12-19) entire text, all drawings	1-7

☐ Further documents are listed in the continuation of Box C.
 ☒ See patent family annex.

* Special categories of cited documents:

“A” document defining the general state of the art which is not considered to be of particular relevance

“E” earlier application or patent but published on or after the international filing date

“L” document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)

“O” document referring to an oral disclosure, use, exhibition or other means

“P” document published prior to the international filing date but later than the priority date claimed

“T” later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention

“X” document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone

“Y” document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art

“&” document member of the same patent family

Date of the actual completion of the international search

19 April 2022

Date of mailing of the international search report

26 April 2022

Name and mailing address of the ISA/JP

Japan Patent Office (ISA/JP)
 3-4-3 Kasumigaseki, Chiyoda-ku, Tokyo 100-8915
 Japan

Authorized officer

Telephone No.

INTERNATIONAL SEARCH REPORT
Information on patent family members

International application No.
PCT/JP2022/005497

Patent document cited in search report			Publication date (day/month/year)	Patent family member(s)	Publication date (day/month/year)
JP	2001-265596	A	28 September 2001	(Family: none)	
US	2019/0384800	A1	19 December 2019	US 2017/0228239 A1 entire text, all drawings	

REFERENCES CITED IN THE DESCRIPTION

This list of references cited by the applicant is for the reader's convenience only. It does not form part of the European patent document. Even though great care has been taken in compiling the references, errors or omissions cannot be excluded and the EPO disclaims all liability in this regard.

Non-patent literature cited in the description

- **YUSUKE KOYANAGI**. Developing a Framework for Individual Causal Discovery and its Application to Real Marketing Data. *The Japanese Society for Artificial Intelligence 18th Special Interest Group on Business Informatics*, March 2021, http://sig-bi.jp/doc/18thSIG-BI2021/18thSIG-BI2021_paper13.pdf> **[0005]**
- **TAKEAKI UNO ; TATSUYA ASAI ; YUZO UCHIDA ; HIROKI ARIMURA**. An Efficient Algorithm for Enumerating Closed Patterns in Transaction Databases. *Discovery Science*, 2004, vol. 3245, 16-31 **[0025]**